

# Multivariate Analysis of Stock Prices

Group 11

Abhijit Krishna Menon

Manya Raman

December 9, 2019

## CONTENTS

1	Background and Introduction	2
1.1	The Problem . . . . .	2
1.2	Goal of the study . . . . .	2
1.3	Possible Solution . . . . .	2
2	Data Exploration and Visualization	3
3	Data Preparation and Pre-processing	5
3.1	Stock Data . . . . .	6
3.2	Company News Data . . . . .	6
3.3	Reddit Data . . . . .	6
4	Data Mining Techniques and Implementation	7
4.1	Sentiment Analysis . . . . .	7
4.2	Linear Regression Techniques . . . . .	7
4.3	Neural Networks . . . . .	9
4.4	Recurrent Neural Network - Long Short Term Memory . . . . .	11
5	Performance Evaluation	14
6	Discussion and Recommendation	16
7	Summary	16
8	Appendix: R Code for use case study	17
8.1	Data Collection and Pre Processing. . . . .	17
8.2	Data Visualization . . . . .	21
8.3	Modelling . . . . .	24
9	References	41

---

<sup>1</sup> IE 7275 Data Mining in Engineering

# 1 BACKGROUND AND INTRODUCTION

## 1.1 The Problem

The stock market is not a clear cut equation and it depends on a lot of factors and not just numbers. In this world dominated by fake news and headlines, a lot of our economic decisions are made through what we see taking place around us. One such example is that negative news will normally cause individuals to sell stocks. A single headline of a company does not only change the stocks of that only but also affects the other companies related to it. What we are trying to achieve via this project is to determine how much does a single news article affect the stock of a particular company stocks.

## 1.2 Goal of the study

The stock market prediction has been an active area of research for quite a while. However, building a model that takes into consideration every factor is still a challenging problem. Apart from historical prices, the current stock market is affected by news articles about the company, general news and many other micro economic and macroeconomic factors. There are several models out there to predict the stocks based on general news elements or historical prices but none have taken into consideration all company specific news or all of these factors together.

In this project, our goal is to analyse how much of a role do these factors play while attempting to predict the next day's stock price. Thus, making a model with checking all possible computations and permutations that improves its accuracy and delivers an efficient model.

## 1.3 Possible Solution

Our project will focus on two tracks that join into one rail at the end.

The first track will be to build a sentiment analysis model using the common news articles and the company specific news articles.

The second track is to utilize the numerical data available to us about the stocks and linearly model this data with our final target.

Finally, using the sentiment generated from the first track we would like to see if the sentiment scores generated actually influence our final output as compared to when we do not utilize it for predictions.

### 1.3.1 Github Repository

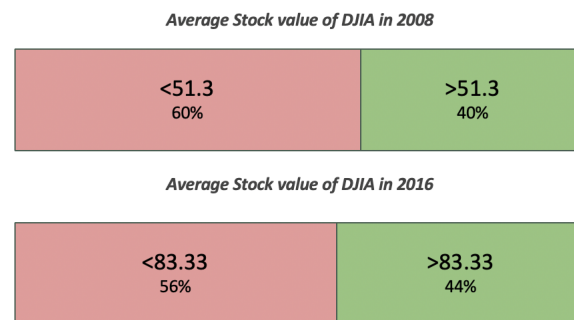
The link to our github repository including the code and other details of our project is <https://github.com/akmenon1996/Multivariate-Stock-Market-Analysis>

## 2 DATA EXPLORATION AND VISUALIZATION

The stock data collected from Yahoo API of DJIA(Dow Jones Industrial Average) companies for 8 years has data missing on weekends. The news data collected from all the resources were filtered only for the business section which was a key element of the API to generate relevant news for our purpose. These news items are not present for every day of the week because of the lack of news articles for a single company across several days by a single news source. In such a case the sentiment for that day is taken as zero.

After cleaning and transforming the data to a proper format, the next step was to discover meaningful relations and summarize the main characteristics of the data set. For this purpose we performed a detailed Exploratory Data Analysis(EDA) on our data.

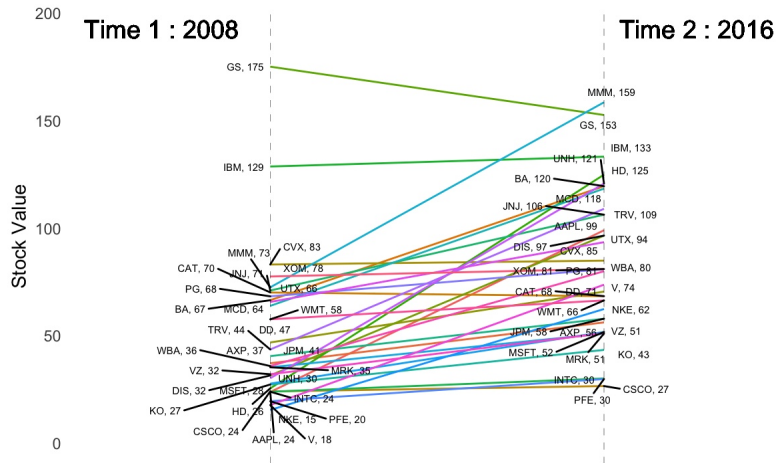
First we analyzed the average stock value in 2008, that was 51.30 and the DJIA companies above average were 40%. But in 2016 the metric changed and the average stock value increased becoming 83.33. But the interesting part we analyzed was that the companies above average are now different than some of those that were in 2008.



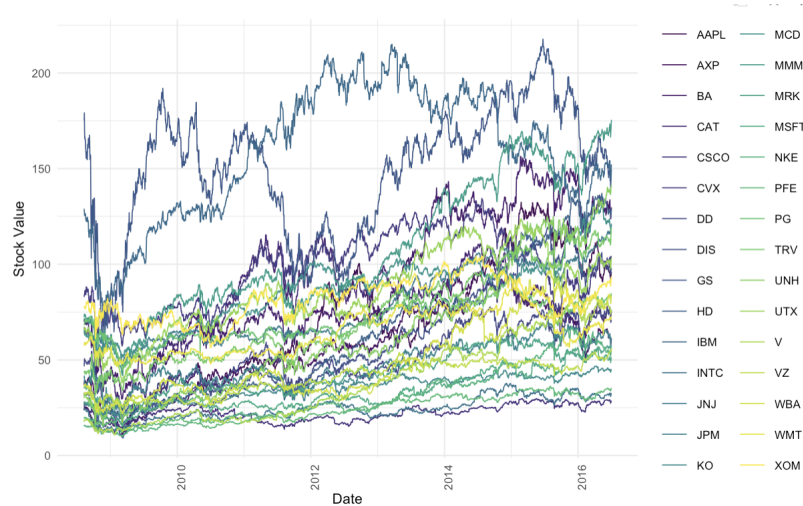
**Figure 1:** Summary of Average Stock Value change in 8 years.

To understand this increase in percentage and change in companies stock above average we did an in-depth analysis country-wise. Therefore, figure 2 and figure 3 shows the change in stock value of 30 companies, from where we can see that companies like Apple, Disney, Home Depot etc. which were below average in 2008 are above average in 2015, and companies like Walmart, Procter & Gamble and Caterpillar Inc. which were above average in 2008 are now below average.

Also, it is true that not all these companies have always been a part of the DJIA. These are only the current DJIA list, however since the DJIA list is ever fluctuating, just these companies is in no way indicative of the entire stock market.

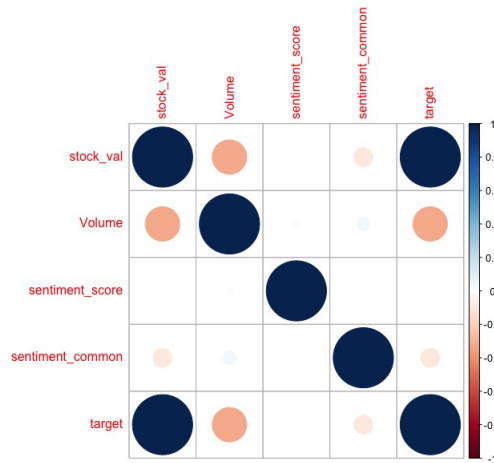


**Figure 2:** Slope Plot for change in stock value over time



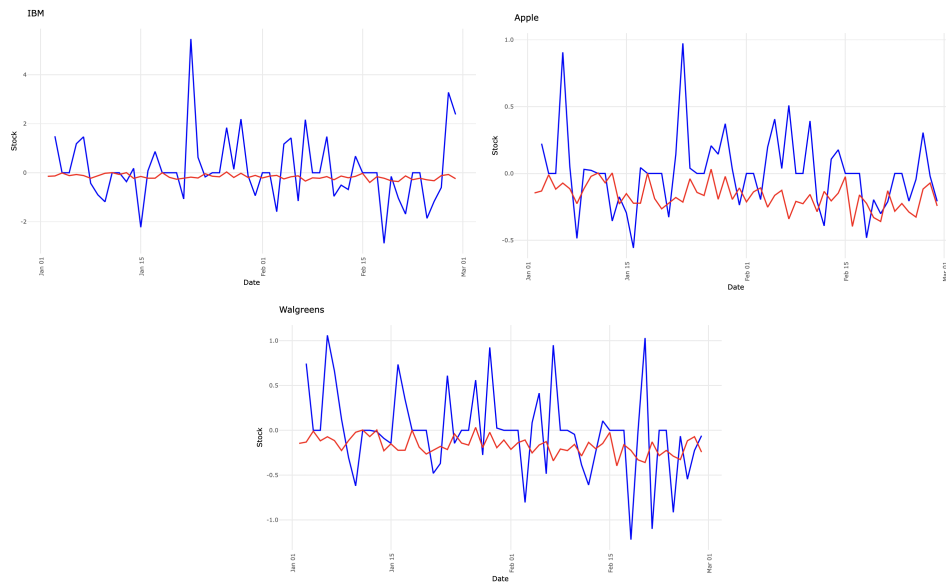
**Figure 3:** Line Plot for change in stock value over time

After the visualization of the stock data we started visualizing the correlation between different variables as seen from the correlation plot in figure 4. We analyzed that there is correlation between stock value and common sentiment thus we did an in-depth analysis.



**Figure 4:** Correlation Plot between input variables and target.

To visualize correlation between stock value and sentiment we first removed the trend from the stock data and lagged the stock value by 1 as today's news sentiment will affect the stock value of tomorrow. Figure 5 displays the zoomed version of Apple, IBM and Walgreens stocks and news sentiments respectively. We can see that from the figures that there is a directly proportional relationship between news sentiment and stock value.



**Figure 5:** Sentiment (Red) and Stock (Blue) for IBM, Apple and Walgreens.

### 3 DATA PREPARATION AND PRE-PROCESSING

The data is collected from 3 different sources and pre-processed in different manner :

### 3.1 Stock Data

We have collected data of 30 companies over a period of 2006 to 2016 from Yahoo finance API. The attributes are: Opening price, Closing price, high and low. The stock data is absent for weekends and other holidays when the market is closed. In order to complete the data, we approximate the missing values using concave function as the stock data usually follows this function. So, if the DJIA value on a given day is  $x$  and the next available data point is  $y$  with  $n$  days missing in between, we approximate the missing data by estimating the first day after  $x$  to be  $(y+x)/2$  and then following the same method recursively until all gaps are filled.

### 3.2 Company News Data

The data for 30 companies under Dow Jones Industrial Average (DJIA) is collected from the NY-Times API. The attributes are : created\_time, snippet, headline, news desk and company name. The problem with data collection using the API was that it allows only 110 articles of any kind for a given date. We have mined data for over 6 years across 30 companies. We achieved this by writing a recursive function in R which takes into consideration the number of hits the API can make in a given minute and the number of pages it should look to as well as the year and company that we are looking at. After getting the data we filtered out the data on the condition of the section of news it originates from. The data was collected for every company individually thus we first joined all the companies files and created one with all. After that, a single string was formed from concatenating all the articles headlines for a single day company-wise as some companies have multiple headlines in a single day. Finally we removed all the punctuation's and byte characters thus getting a final cleaned and formatted version of company-wise news data.

### 3.3 Reddit Data

The data is collected from Kaggle. This data contains historical news headlines from Reddit World-News Channel. They are ranked by reddit users' votes, and only the top 25 headlines are considered for a single date. The first column is the "date", and the second column is the "news headlines". Hence, there are 25 lines for each date. For pre-processing, a single string was formed from concatenating all the articles headlines for a single day.

After all three files are collected they are combined into one to do all the visualizations and data mining. The news data was merged with Dow Jones Industrial Average (DJIA) companies stock index value on appropriate date (lag by 1 day) to get final data frame.

---

## 4 DATA MINING TECHNIQUES AND IMPLEMENTATION

### 4.1 Sentiment Analysis

For the calculation of the polarity score of the two news headlines of each company we used an in-built sentiment score calculator that came as a part of the 'sentimentr' package on R. We pivoted our data so as to have two columns, one for the sentiment of the common news and one for the company specific news. We then calculated the sentiment for each of these news items and added them to the respective columns by day for each company. We had the sentiment for every day of the common news, however the company specific news was mostly 0s because not every day had a news item for every company, also we could not collect a lot of news data across multiple sources. If the sentiment for that day is NA then we decided to set the sentiment to 0 which essentially means the polarity for that day is neutral. Our assumption here is that a neutral polarity news items do not affect the stock price at all.

### 4.2 Linear Regression Techniques

For building the linear models, we first plotted scatter plots amongst all the variables with the target variable to visualize linearity(if any). From the scatter plots we notice that the only variable with a linear relation to the target value is the previous day's stock value, which was fairly expected. The volume of stocks traded has an exponential relation to the target after performing a log transformation and the news data does not have any heavy correlation. The scatter plots can be seen as below.

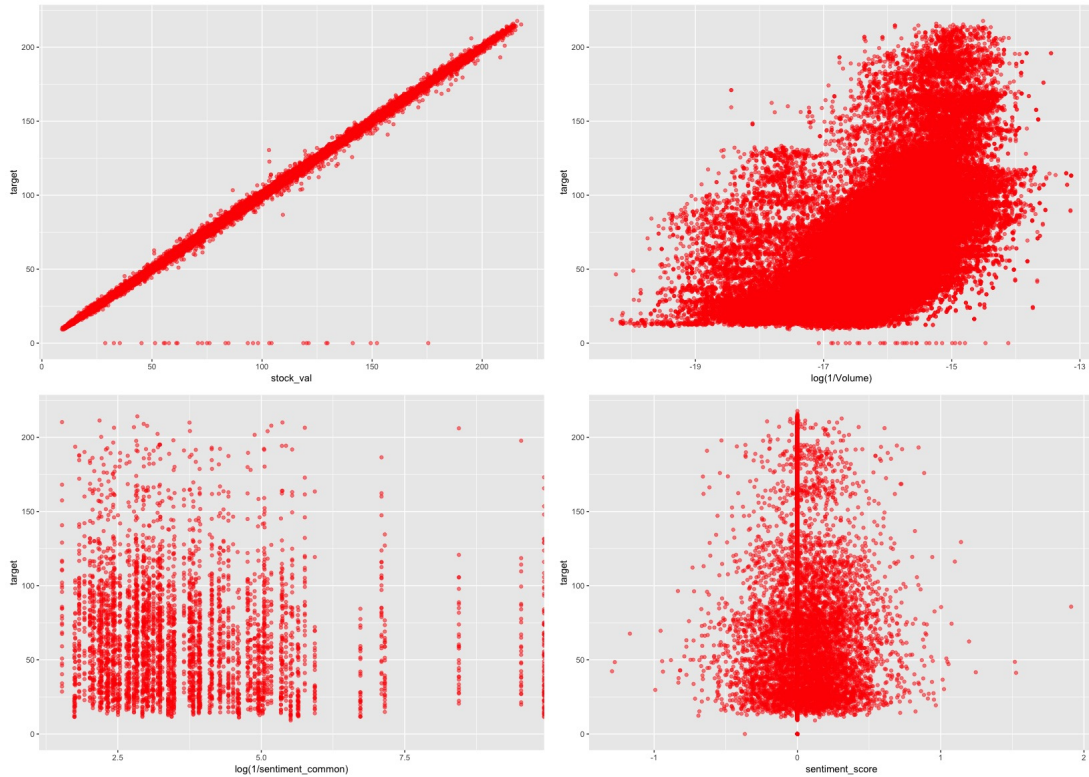
#### 4.2.1 Generalised Linear Model

The formula that we used to build the initial model on R was the following.

```

1
2 Call:
3 lm(formula = target ~ stock_val + log(1/Volume) + sentiment_score +
4     sentiment_common, data = train_stock_df)
5
6 Coefficients:
7     (Intercept)      stock_val      log(1/Volume)  sentiment_score
8     sentiment_common
9     0.271537      0.999769      0.014899      0.116579
10    0.002159
11
12 summary(linear_model)
13
14 Call:
15 lm(formula = target ~ stock_val + log(1/Volume) + sentiment_score +
16     sentiment_common, data = train_stock_df)

```



**Figure 6:** Scatter plots between input variables and target variable.

```

16
17 Residuals:
18      Min       1Q   Median       3Q      Max
19 -22.7602  -0.1377  -0.0146   0.1661  27.3592
20
21 Coefficients:
22             Estimate Std. Error t value Pr(>|t|)
23 (Intercept)   0.2715367   0.0643611    4.219 2.46e-05 ***
24 stock_val     0.9997692   0.0001047  9549.984 < 2e-16 ***
25 log(1/Volume) 0.0148994   0.0037255    3.999 6.36e-05 ***
26 sentiment_score 0.1165795   0.0434593    2.682 0.00731 **
27 sentiment_common 0.0021588   0.0228402    0.095 0.92470
28
29 Signif. codes:  0  '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1
30
31 Residual standard error: 0.7771 on 59155 degrees of freedom
32 Multiple R-squared:  0.9995, Adjusted R-squared:  0.9995
33 F-statistic: 3.082e+07 on 4 and 59155 DF, p-value: < 2.2e-16

```

From the summary of the linear model we see that the most important variables are:

- Stock Value
- 1/Volume



- Sentiment Score

On testing and generating our accuracy, we get a shocking correlation accuracy of 99%, however this is obviously not the right way to model the Time series data that we have with us. The reason being the high amount of trend, seasonality, and the high correlation between the he previous day's stock to the next day's affects the accuracy of the output. This is thus not the right model to be using for our data.

	actuals	predicted
actuals	1.0000000	0.9963687
predicted	0.9963687	1.0000000

#### 4.2.2 Support Vector Regressor

```

1 SVR_model <- svm(target ~ stock_val+log(1/Volume)+sentiment_score+sentiment_
  common, data=train_stock_df) # build linear regression model on full data
2 print(SVR_model)
3
4 Call:
5 svm(formula = target ~ stock_val + log(1/Volume) + sentiment_score + sentiment
  _common, data = train_stock_df)
6
7
8 Parameters:
9   SVM-Type:  eps-regression
10  SVM-Kernel: radial
11      cost:   1
12   gamma:   0.25
13  epsilon:   0.1
14
15
16 Number of Support Vectors: 810

```

Simillarly with the SVR we get a hiogh correlation accuracy of 99%. Again because of the earlier reasons, it is not the right model to use for this sort of data.

	actuals	predicted
actuals	1.0000000	0.9929886
predicted	0.9929886	1.0000000

### 4.3 Neural Networks

Since the linear models were clearly not the way to go ahead with the data that we had with us, we decided to push our data through a bunch of neural networks and test the accuracy of the data predicted and the fit with the model.

For our experimentation we built different models with adding more variables at a time to see the difference in the prediction accuracy. This way, we can see the effect(if any) of the news variables on the final outcome.

The model used across all iterations of model bulding were constant. Just the variables inserted into the models variable.

For the base case, we tested 3 dense layers with the rectified linear unit(RELU) activation function. The loss was measured by the mean square error and the metrics we used to check accuracy is the mean absolute error.

```

1      build_model <- function() {
2
3      model <- keras_model_sequential() %>%
4        layer_dense(units = 64, activation = "relu",
5                      input_shape = dim(train_data)[2]) %>%
6        layer_dense(units = 64, activation = "relu") %>%
7        layer_dense(units = 1)
8
9      model %>% compile(
10        loss = "mse",
11        optimizer = optimizer_rmsprop(),
12        metrics = list("mean_absolute_error")
13      )
14
15      model
16    }
17
18    model <- build_model()
19    model %>% summary()
20
21
22
23    Model: "sequential_7"
24
25 -----
26
27 Layer (type)                                     Param #      Output
28 -----
29 dense_21 (Dense)                                256          (None ,
30   64)
31
32 dense_22 (Dense)                                4160         (None ,
33   64)
34
35 dense_23 (Dense)                                65          (None ,
36   1)
37
38 Total params: 4,481
39 Trainable params: 4,481

```

35 Non-trainable params: 0

36

The accuracy of these sections will be discussed in the Performance and Evaluation Section.

#### 4.4 Recurrent Neural Network - Long Short Term Memory

From the results of the linear models and the traditional neural network we realise that the other variables are just not being captured enough by the model, because of the stupendously high correlation between the previous day's data and the next day's data which is what was our target.

Thus it was necessary to build a different model, while still not going completely down the track of time series modelling. We decided to build an LSTM model, which is a kind of recurrent neural network.

What is a RNN or a recurrent Neural network? A brief explanation from one of our sources might be able to explain this better.

---

##### Recurrent Neural Networks

Humans don't start their thinking from scratch every second. As you read this essay, you understand each word based on your understanding of previous words. You don't throw everything away and start thinking from scratch again. Your thoughts have persistence.

Traditional neural networks can't do this, and it seems like a major shortcoming. For example, imagine you want to classify what kind of event is happening at every point in a movie. It's unclear how a traditional neural network could use its reasoning about previous events in the film to inform later ones.

Recurrent neural networks address this issue. They are networks with loops in them, allowing information to persist.

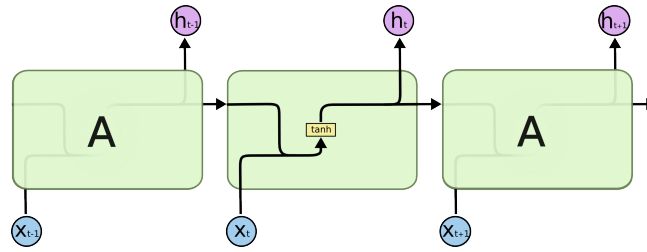
---

Now that we have an understanding of what an RNN is, we get a clearer picture of what role it plays for our data. Although we are trying to generate a linear relationship between our data points, it simply does not exist. This is why the previous data points are as important as the current one while making a decision.

The LSTM network allows us "remember" previous data points as a history and add that as an input variable while generating the output of the present data set.

```

1         rmsprop = optimizer_rmsprop(lr=0.00001, rho=0.9, epsilon=1e-08)
2 model <- keras_model_sequential()
3
4 model %>%
5     layer_lstm(units           = 100,
6               input_shape     = c(dim(train_X)[2], dim(train_X)[3]),
7               return_sequences = TRUE,
```



**Figure 7:** Network diagram of the LSTM model.

```

8         ) %>%
9         layer_dropout(0.5)%>%
10        layer_lstm(units = 75,
11                  return_sequences = FALSE,
12                  ) %>%
13        layer_dropout(0.5)%>%
14        layer_flatten()%>%
15        layer_dense(units = 64, activation = "relu") %>%
16        layer_dropout(0.5)%>%
17        layer_dense(units = 1)
18
19 model %>%
20   compile(loss = 'mae', optimizer = rmsprop)
21
22 model

```

```

25 Model
26 Model: "sequential_3"

```

Layer (type)	Shape	Param #	Output
lstm_6 (LSTM)	1, 100	42000	(None ,
dropout_9 (Dropout)	1, 100	0	(None ,
lstm_7 (LSTM)	75)	52800	(None ,
dropout_10 (Dropout)	75)	0	(None ,
flatten_3 (Flatten)	75)	0	(None ,

```

39 dense_6 (Dense)                                     (None ,
    64)                                                4864
40 -----
41 dropout_11 (Dropout)                               (None ,
    64)                                                0
42 -----
43 dense_7 (Dense)                                     (None ,
    1)                                                65
44 =====
45 Total params: 99,729
46 Trainable params: 99,729
47 Non-trainable params: 0
48 -----

```

For building the optimum model we used 2 LSTM layers, separated by a drop out layer to avoid overfitting. Followed a flatten layer, a dense layer and another dense layer of 1 that gives us the final output.

The input to the LSTM models is slightly different. The input is fed in as a 3d matrix. With the first dimensions being total number of records, the second being the numebr of time steps to be remembered and the 3rd dimensions as the input variables.

We trained this model over 50 epochs and a batch size of 10.

```

1 epochs = 50
2
3
4 history <- model %>% fit(
5   train_X ,
6   train_y ,
7   epochs = epochs ,
8   batch_size = 10,
9   validation_split = 0.2 ,
10  verbose = 1,
11  shuffle = FALSE
12 )

```

Using the LSTM, gave us satisfaction that the model is not over fit and is not overly dependant on a single variable. Using this model we made a few predictions, which we will cover over the next section.

## 5 PERFORMANCE EVALUATION

The model was split into 2 parts, the training data as well as the validation data in a 70-30 split. The performance evaluation of our various models are as follows:

- Linear Regression

Mean Absolute Error of 4.82. Correlation Accuracy between Actual and Predicted values was 0.996.

- Support Vector Regression

Mean Absolute Error of 3.22. Correlation Accuracy between Actual and Predicted values was 0.992.

- Neural Network - Only Stock Value and Volume.

Mean absolute Error of 1.28.

- Neural Network - Stock Value, Volume and Sentiment of common news data.

Mean absolute Error of 1.56.

- Neural Network - Only Sentiment columns

Mean Absolute Error of 5.78

- Neural Network - All columns

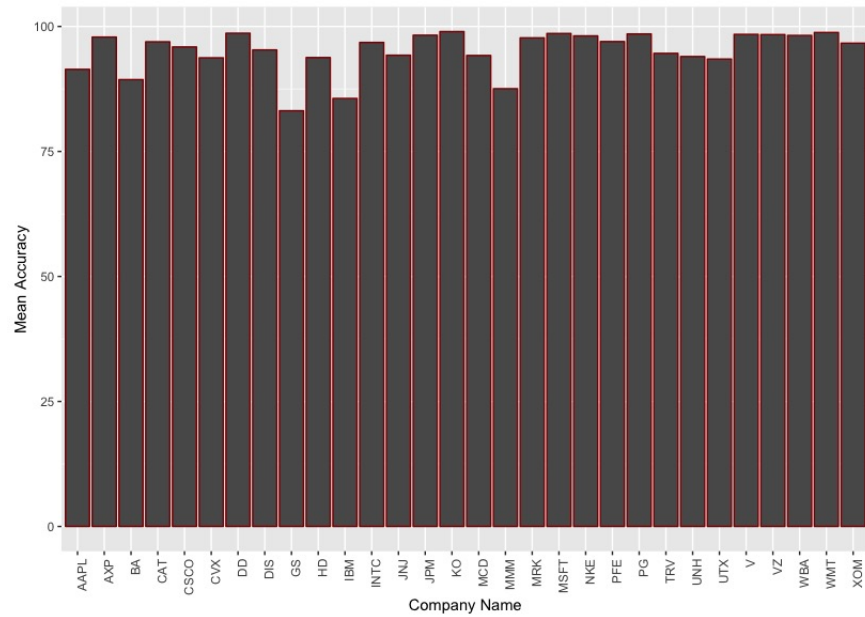
Mean Absolute Error of 1.0002.

- Recurrent Neural Network - LSTM (All columns)

Mean Absolute Error of 0.45

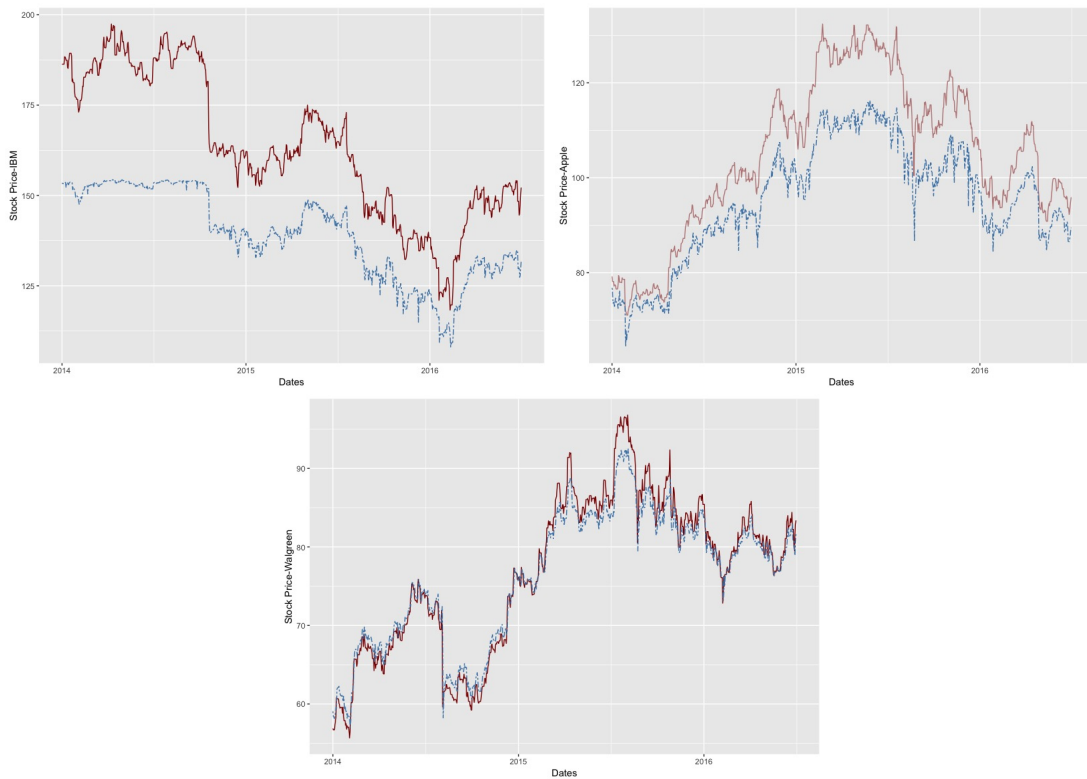
Thus we can see that the LSTM model proves to be the most accurate in terms of the MAE metric.

The company wise mean accuracy is also depicted below.



**Figure 8:** Company wise accuracy

The results of a few companies predictions are below. The red lines are the predicted values.



**Figure 9:** Predictions of data for IBM, Apple and Walgreens.

## 6 DISCUSSION AND RECOMMENDATION

Our primary issue with the data was that it is a time series model that we are trying to evaluate linearly. To get a better performance and a better understanding of the data, we need to first work on a few time series techniques of removing trend, seasonality etc and work with a few time series models such as ARIMA etc. Next, we notice that from our data, the company news sentiment does not seem to be playing a massive role in the stock prices, however we know that, that is not true. Company news does indeed play a massive role in the stock prices. The problem with our data was that we were not able to collect too much of company news data for it to be viable from the New York times API. Solving this issue should have us gain a larger understanding on the effect that the news has on the stock prices.

---

## 7 SUMMARY

From our analysis, we see that sentiment does play a slight role on the price of the next day's stock, although not a whole lot. This is mainly because our news data set is not the most exhaustive resource there is. There is scope for a lot more improvement on the analysis of the affect of news sentiment on the stock prices.

Amongst the models we notice that the LSTM model performs considerable better than the other models. This does not come as a surprise since the LSTM model is the only model that takes into consideration the time series nature of this data set.

---



## 8 APPENDIX: R CODE FOR USE CASE STUDY

Our whole code can be found via the GIT HUB repository below. Some of the major chunks have been added below.

### 8.1 Data Collection and Pre Processing.

#### 8.1.1 *Pre-processing of News data from NYTimes API*

```

1  title: "R Notebook"
2  output: html_notebook
3
4  ---
5  { r }
6  library(dplyr)
7
8
9  { r }
10 library('googledrive')
11 api = "Nql4mlbzy44BETXG70DGesoOoXlnXeKH"
12
13
14
15
16 { r }
17 if (!require("jsonlite")) install.packages("jsonlite")
18 library(jsonlite)
19 #####
20 #####      function — search news article with API
21 #####
22 nytime = function (keyword, year) {
23   searchQ = URLEncode(keyword)
24   url = paste('http://api.nytimes.com/svc/search/v2/articlesearch.json?q=',
25               searchQ,
26               '&begin_date=', year, '0101&end_date=', year, '1231&api-key=', api,
27               sep="")
28   #get the total number of search results
29   initialsearch = fromJSON(url, flatten = T)
30   maxPages = round((initialsearch$response$meta$hits / 10)-1)
31
32   #try with the max page limit at 10
33   maxPages = ifelse(maxPages >= 10, 10, maxPages)
34   #creat a empty data frame
35   df = data.frame(id=as.numeric(), created_time=character(), snippet=character(),
36                   ,
37                   headline=character(), news_desk = character(), company_name =
38                   character())
39   #save search results into data frame

```

```

36 r <- NULL
37 attempt <- 1
38 while( is.null(r) && attempt <= 3 ) {
39   print(paste0("Try:", attempt))
40   attempt <- attempt + 1
41   Sys.sleep(10)
42   try(
43     for(i in 0:maxPages){
44       #get the search results of each page
45       nytSearch = fromJSON(paste0(url, "&page=", i), flatten = T)
46       temp = data.frame(id=1:nrow(nytSearch$response$docs),
47                         created_time = nytSearch$response$docs$pub_date,
48                         snippet = nytSearch$response$docs$snippet,
49                         headline = nytSearch$response$docs$headline.main,
50                         news_desk = nytSearch$response$docs$news_desk,
51                         company_name = keyword)
52       df=rbind(df, temp)
53       Sys.sleep(10)
54     }
55   )
56 }
57 return(df)
58 }
59 ```
60
61 ```{r}
62 company_list = c('Microsoft', 'Nike', 'Pfizer', 'Procter & Gamble', 'The
63                 Travelers Companies', 'United Health Group', 'United Technologies', 'Verizon',
64                 'Visa Inc.', 'Walmart', 'Walgreens Boots Alliance', 'The Walt Disney Company'
65                 ')
66
67 years = c(2008,2009,2010,2011,2012,2013,2014,2015)
68 ```
69
70 ```{r}
71 companies_df = data.frame(id=as.numeric(), created_time=character(), snippet=
72                           character(),
73                           headline=character(), news_desk = character(), company_name =
74                           character())
75 for (company in company_list){
76   cat(paste0("Currently getting the data for ", company, "\n"), file = "Company_
77     Data/data_gen.txt", append = TRUE)
78   print(paste0("Currently getting the data for ", company))
79   company_df = data.frame(id=as.numeric(), created_time=character(), snippet=
80                           character(),
81                           headline=character(), news_desk = character(), company_name =
82                           character())
83   for(year in years){
84     print(paste0(" ", year))

```

```

79   cat(paste0("      Currently getting the data for the year ",year,"\n"),file
    = "Company_Data/data_gen.txt",append = TRUE)
80   df = nyttime(company,year)
81   company_df <- rbind(company_df,df)
82 }
83 file1 = paste0("Company_Data/",company,".csv")
84 print(file1)
85 write.csv(company_df,file = file1)
86 companies_df <- rbind(companies_df,company_df)
87 cat(paste0("Succesfully got the data for ",company,"\n"),file = "Company_
    Data/data_gen.txt",append = TRUE)
88 }
89 write.csv(companies_df,file = "Company_Data/Final_Dataframe.csv")
90
91 ""
92

```

### 8.1.2 News data Preprocessing

```

1  ———
2  title: "News data Preprocessing"
3  output: html_document
4  ———
5
6  ```{r setup, include=FALSE}
7  library(readr)
8  library(plyr)
9  library(dplyr)
10 ```
11
12 ## R Markdown
13
14 #Creating combined dataset by including all companies news data. NYTimes Data:
15
16 ```{r concatenating_files}
17
18 file_list <- list.files()
19 for (file in file_list){
20
21   # if the merged dataset doesn't exist, create it
22   if (!exists("dataset")){
23     dataset <- read.table(file, header=TRUE, sep=",")
24   }
25
26   # if the merged dataset does exist, append to it
27   if (exists("dataset")){
28     temp_dataset <-read.table(file, header=TRUE, sep=",")
29     dataset<-rbind(dataset, temp_dataset)
30     rm(temp_dataset)
31   }

```

```

32 }
33 }
34 final_df <- dataset
35 rm(dataset)
36
37 ““
38
39
40 ““{r preprocessing data}
41
42 final_df_new<- final_df
43 final_df_new$created_time<-as.Date( final_df$created_time ,format="%Y-%m-%d")
44
45 company_data <-final_df_new %>%
46     filter(news_desk == 'Business' & created_time >= '2008-08-08
47     ' & created_time <= '2016-07-01')%>%
48     select(company_name,created_time ,headline)
49
50 company_data<-ddply(company_data , .(company_name,created_time) , summarise ,
51     headline=paste0(headline , collapse="; "))
52
53 levels (company_data$company_name)
54
55 write.csv(company_data , file = "/Users/manyaraman/Desktop/Comapny_Data.csv" ,
56     row.names = FALSE)
57 ““
58
59 #Preprocessing and formatting Reddit data :
60
61 ““{r reddit_news_data}
62
63 Combined_news<-read_csv( 'Combined_News_DJIA.csv' )
64 Combined_news<-Combined_news[ , -2]
65
66 for(i in colnames(Combined_news[ , -1])){
67     Combined_news[[ i ]]<-str_sub(Combined_news[[ i ]] , 2)
68 }
69
70 Combined_news_final <- unite(Combined_news, com_headlines , Top1:Top25 , sep = "
71     ;" , remove = TRUE)
72
73 Reddit_news<-read_csv( 'RedditNews.csv' )
74 Reddit_news<-ddply(Reddit_news , .(Date) , summarise , News=paste0(News , collapse
75     ="; "))
76
77 Reddit_news$News<-str_replace_all(Reddit_news$News,c("b'" = "" , "b\" = ""))
78
79 write.csv(Reddit_news , file = "/Users/manyaraman/Desktop/Combined_news_final_
80     reddit_wknd.csv" , row.names = FALSE)
81 ““

```

77

## 8.2 Data Visualization

```

1  ———
2  ———
3  title: "Project-visualization"
4  output: html_document
5  ———
6
7  ““{r setup, include=FALSE}
8  knitr::opts_chunk$set(echo = TRUE)
9  library(ggplot2)
10 library(dplyr)
11 library(plyr)
12 library(readxl)
13 library(plotly)
14 library(tidyverse)
15 library(ggrepel)
16 library(gganimate)
17 library(tseries)
18 library(forecast)
19 ““
20
21 #Line plot of all companies stocks:
22 ““{r pressure, warning = false,echo=FALSE}
23
24 stock_plot <- ggplot(
25   final_df,
26   aes(Date, stock_val, group = company_name, color = factor(company_name))
27 ) +
28   geom_line(size=0.3) +
29   scale_color_viridis_d() +
30   labs(x = "Date", y = "Stock Value") +
31   theme(legend.position = "top")+
32   theme_minimal(base_size = 8)+ theme(axis.text.x = element_text(angle = 90))
33
34
35 ggplotly(stock_plot + scale_color_discrete(name="Company"))
36
37 ##stock_plot + geom_point() + transition_reveal(Date)
38
39 ““
40
41 # Stock and Sentiment Analysis Line Plot :
42
43 ““{r Line Plot}
44 final<- final_df%>%filter(company_name == 'AAPL' & Date>='2009-01-01' & Date <
  '2009-03-01')

```

```

45
46 # normalize <- function(x) {
47 # num <- x - min(x)
48 # denom <- diff(range(x))
49 # return (num/denom)
50 # }
51
52 #final$stock_val<- normalize(final$stock_val)
53 #final$sentiment_common<-normalize(final$sentiment_common)
54
55 final_f<- final[-1,]
56 final_f$stock_val<- diff(final$stock_val)
57 final_f$stock_lag<- lag(final_f$stock_val,n=1)
58
59 p <- ggplot() +
60   geom_line(data = final_f, aes(x = Date, y = stock_lag), color = "blue") +
61   geom_line(data = final_f, aes(x = Date, y = sentiment_common), color = "red"
62   ) +
63   scale_color_viridis_d() +
64   labs(x = "Date", y = "Stock", title = "Walgreens") +
65   theme(legend.position = "top")+
66   theme_minimal(base_size = 8)+ theme(axis.text.x = element_text(angle = 90))
67
68 ggplotly(p)
69
70 “““
71 # Companies stock value change over time using sliding bars
72 ““{r Sliding bar plot }
73
74 plotdata <- final_df %>%
75   group_by(Date) %>%
76   mutate(ordering = rank(stock_val)) %>%
77   ungroup()
78
79 p<-ggplot(plotdata ,
80   aes(ordering , group = company_name, color=factor(company_name), fill=
81   factor(company_name), show.legend = FALSE)) +
82   geom_tile(aes(y = stock_val/2,
83   height = stock_val,
84   width = 0.9), alpha = 0.4) +
85   # text on top of bars
86   geom_text(aes(y = stock_val, label = as.integer(stock_val)), hjust = -0.2) +
87   # geom_text(aes(y = 0, label = country, hjust = 3)) +
88   # text in x-axis (requires clip = "off" in coord_cartesian)
89   geom_text(aes(y = 0, label = company_name, color="black"), hjust = 0, show.
90   legend = FALSE) +
91   coord_flip(clip = "off", expand = TRUE) +
92   enter_fade() +
93   exit_shrink() +
94   coord_flip()+

```

```

93   scale_color_viridis_d(name="", guide=FALSE)+
94   scale_fill_viridis_d(name="",guide=FALSE)+
95   ylim(0, 200) +
96   # theme_tufte(14,"Avenir")+
97   theme_classic() +
98   # guides(color=F, fill=F)+
99   labs(title = "Year: {closest_state}", y="Stock", x="Company" ) +
100  theme(plot.title = element_text(hjust = 0.5, size = 24),
101        axis.ticks.y = element_blank(),
102        axis.text.y = element_blank()) +
103  transition_states(states = Date, transition_length = 2, state_length = 1) +
104  # transition_time(year)+
105  ease_aes('cubic-in-out')
106
107  animate(p, nframes = 160, fps = 20, end_pause = 20, width = 500, height = 900)
108  #use anim_save(filename) to save
109
110  anim_save("animation_le4.gif", animation = last_animation())
111
112  ‘‘‘
113  # Slope plot for change in stock value over time
114
115  ‘‘{r Slope plot}
116  fc<-final_df %>% filter(Date == '2008-08-08' | Date == '2016-03-01')%>%dplyr::
117    select(Date, stock_val, company_name)
118  fcs<-spread(fc, key = Date, value = stock_val)
119
120  left_label <- paste(fcs$company_name, round(fcs$'2008-08-08'), sep=", ")
121  right_label <- paste(fcs$company_name, round(fcs$'2016-03-01'), sep=", ")
122
123  p <- ggplot(fcs) + geom_segment(aes(x=1, xend=2, y='2008-08-08', yend
124    ='2016-03-01', color=factor(company_name)), size=.5, show.legend=F) +
125    geom_vline(xintercept=1, linetype="dashed", size=.1) +
126    geom_vline(xintercept=2, linetype="dashed", size=.1) +
127    labs(x="", y="Stock Value") + # Axis labels
128    xlim(.5, 2.5) + ylim(0,(1.1*(max(fcs$'2008-08-08', fcs$
129    '2016-03-01'))))
130
131  p <- p + geom_text_repel(label=left_label, y=fcs$'2008-08-08', x=rep(1, NROW(
132    fcs)), hjust=0.1, size=2)
133  p <- p + geom_text_repel(label=right_label, y=fcs$'2016-03-01', x=rep(2, NROW(
134    fcs)), hjust=0.1, size=2)
135  p <- p + geom_text(label="Time 1 : 2008", x=1, y=1.1*(max(fcs$'2008-08-08',
136    fcs$'2016-03-01')), hjust=1.2, size=5) # title
137  p <- p + geom_text(label="Time 2 : 2016", x=2, y=1.1*(max(fcs$'2008-08-08',
138    fcs$'2016-03-01')), hjust=-0.1, size=5) # title
139
140  # Minify theme
141  p + theme(panel.background = element_blank(),
142            panel.grid = element_blank(),

```

```

137     axis.ticks = element_blank(),
138     axis.text.x = element_blank(),
139     panel.border = element_blank())
140
141   “ “ “

```

---

## 8.3 Modelling

### 8.3.1 Sentiment Analysis

```

1  ———
2  title: "Sentiment Analyser"
3  output: html_notebook
4  ———
5
6  “ “ { r }
7  library(sentimentr)
8  library(dplyr)
9  “ “ “
10
11 “ “ { r load data }
12 news_data <- read.csv("Merge_Comapny_Data.csv")
13 head(news_data)
14 “ “ “
15
16 “ “ { r }
17 reddit_data <- read.csv('Combined_news_final_reddit_wknd.csv')
18 head(reddit_data)
19 “ “ “
20
21 “ “ { r get_sentiment }
22 news_data = mutate(news_data, news_data_new = sentiment_by(as.character(
23   snippet))$ave_sentiment)
24 “ “ “
25
26 “ “ { r }
27 news_data
28 “ “ “
29
30 “ “ { r }
31 d <- density(news_data$news_data_new)
32 plot(d, main="Dist")
33 polygon(d, col="red", border="blue")
34 “ “ “
35
36 “ “ { r }
37 stock_numbers <- read.csv("combined_dataframe_djia.csv")

```



```

37 head(stock_numbers)
38 ""
39
40 ""{ r}
41 stock_numbers_ave <- stock_numbers%>%mutate(average_price = (High+Low)/2)%>%
  select(Date, company_name, High, Low, Volume, average_price)
42 stock_numbers_ave
43 ""
44 ""{ r}
45 # library(Hmisc)
46 # stock_numbers_ave$lagged <- Lag(stock_numbers_ave$average_price, +1)
47 # stock_numbers_ave
48
49 stock_numbers_ave <- stock_numbers_ave %>%group_by(company_name) %>%mutate(
  target = dplyr::lead(average_price, n = 1, default = NA))%>%ungroup()
50 ""
51
52
53
54 ""{ r}
55 news_data
56 ""
57
58 ""{ r}
59 levels(news_data$company_name)
60 ""
61
62 ""{ r}
63 levels(stock_numbers_ave$company_name)
64 ""
65
66 ""{ r}
67 company_ticker_dict = list("AAPL"="Apple",
68                             "AXP"="American Express",
69                             "BA"="Boeing",
70                             "CAT"="Caterpillar Inc.",
71                             "CSCO"="Cisco",
72                             "CVX"="Chevron Corporation",
73                             "DD"="Dow",
74                             "DIS"="The Walt Disney Company",
75                             "GS"="Goldman Sachs",
76                             "HD"="The Home Depot",
77                             "IBM"="IBM",
78                             "INTC"="Intel",
79                             "JNJ"="Johnson & Johnson",
80                             "JPM"="JPMorgan Chase",
81                             "KO"="Coca-Cola",
82                             "MCD"="McDonald's",
83                             "MMM"="3M",
84                             "MRK"="Merck & co",
85                             "MSFT"="Microsoft",

```

```

86         "NKE"="Nike" ,
87         "PFE"="Pfizer" ,
88         "PG"="Procter & Gamble" ,
89         "TRV"="The Travelers Companies" ,
90         "UNH"="United Health Group" ,
91         "UTX"="United Technologies" ,
92         "V"="Visa Inc." ,
93         "VZ"="Verizon" ,
94         "WBA"="Walgreens Boots Alliance" ,
95         "WMT"="Walmart" ,
96         "XOM"="ExxonMobil" )
97
98 #stock_numbers_ave%>%mutate(company_full_name = company_ticker_dict[paste0
99   ("'",company_name,"'") ])
100
101 company_name_list = list()
102 for (i in stock_numbers_ave$company_name){
103   company_name_list <- append(company_name_list,company_ticker_dict[i])
104 }
105
106 “““
107
108 ““{ r}
109 company_names_list <- stack(company_name_list)$values
110 “““
111
112 ““{ r}
113 stock_numbers_ave <- stock_numbers_ave%>%mutate(company_name_full = company_
114   names_list)
115 head(stock_numbers_ave)
116
117 ““{ r}
118 news_data%>%filter(company_name == 'IBM')
119 “““
120
121 ““{ r}
122 a <- left_join(stock_numbers_ave, news_data, by = c("Date" = "created_time",
123   "company_name_full" = "company_name"))
124 a
125
126 ““{ r}
127 test <- left_join(a, reddit_data, by = c("Date" = "Date"))
128 “““
129
130 ““{ r}
131 test_sentiment = mutate(test, sentiment_common = sentiment_by(as.character(
132   News))$ave_sentiment)
133 test_sentiment
134 “““

```

```

133
134 ““{ r }
135 write.csv(test_sentiment,"final_sentiment_data_with_reddit.csv")
136 ““
137
138
139
140 ““{ r }
141 #final_df <- na.omit(a)
142 final_df <- test_sentiment %>% mutate(sentiment_score = news_data_new) %>% select(
    Date, company_name, company_name_full, High, Low, Volume, sentiment_score,
    sentiment_common, target)
143 ““
144
145 ““{ r }
146 final_df <- read.csv('Final_data_with_reddit.csv')
147 ““
148
149 ““{ r }
150 write.csv(final_df, 'Final_data_with_reddit.csv')
151 ““
152
153 ““{ r }
154 final_df_non <- final_df
155 final_df_non[is.na(final_df_non)] <- 0
156 ““
157
158 ““{ r }
159 final_df <- na.omit(final_df)
160 ““
161
162 ““{ r }
163 write.csv(final_df_non, 'Final_data_with_reddit.csv')
164 write.csv(final_df, "Final_data_na_filtered")
165 ““

```

### 8.3.2 Linear Modelling

```

1 ———
2 title: "R Notebook"
3 output: html_notebook
4 ———
5 ““{ r }
6 library(keras)
7 library(dplyr)
8 ““
9
10 ““{ r loading_data }
11 stock_df <- read.csv("Final_data_with_reddit.csv")
12 stock_df <- stock_df[,3:11]

```

```

13 head(stock_df)
14 ""
15
16
17 ""{ r}
18 stock_df <- stock_df%>%mutate(stock_val = (High+Low)/2)
19 ""
20
21
22 ""{ r}
23 stock_df <- stock_df%>%select(Date, company_name, company_name_full, stock_val,
    Volume, sentiment_score, sentiment_common, target)
24
25 stock_df
26 ""
27
28 ""{ r}
29 stock_df$Date <- as.Date(stock_df$Date)
30 ""
31
32 ""{ r}
33 sample_size = floor(0.8*nrow(stock_df))
34 set.seed(777)
35 train_stock_df = stock_df%>%filter(Date<'2014-01-01')
36 test_stock_df = stock_df%>%filter(Date>='2014-01-01')
37 ""
38
39 ""{ r}
40 train_stock <- train_stock_df%>%select(stock_val, Volume, sentiment_score,
    sentiment_common)
41 test_stock <- test_stock_df%>%select(stock_val, Volume, sentiment_score,
    sentiment_common)
42 train_target <- train_stock_df%>%select(target)
43 test_target <- test_stock_df%>%select(target)
44 ""
45
46 ""{ r}
47 train_data <- scale(train_stock)
48 col_means_train <- attr(train_data, "scaled:center")
49 col_stddevs_train <- attr(train_data, "scaled:scale")
50 test_data <- scale(test_stock, center = col_means_train, scale = col_stddevs_
    train)
51 ""
52
53 ""{ r}
54 train_data <- as.data.frame(train_stock)
55 test_data <- as.data.frame(test_stock)
56 ""
57
58 ""{ r}
59

```

```

60 plotting_df <- stock_df%>%select(Date,stock_val,Volume,sentiment_score ,
    sentiment_common,target)
61 ""
62
63 ""{ r}
64 plotting_df%>%ggplot() +
65   geom_point(aes(x = stock_val, y = target), color = "red",alpha = 0.5)
66 ""
67
68 ""{ r}
69 plotting_df%>%ggplot() +
70   geom_point(aes(x = sentiment_score, y = target), color = "red",alpha = 0.5)
71 ""
72
73 ""{ r}
74 plotting_df%>%ggplot() +
75   geom_point(aes(x = log(1/sentiment_common), y = target), color = "red",alpha
    = 0.5)
76 ""
77
78
79
80
81 ""{ r}
82 plotting_df%>%ggplot() +
83   geom_point(aes(x = log(1/Volume), y = target), color = "red",alpha = 0.5)
84 ""
85
86
87
88
89 ""{ r}
90 plot(density(plotting_df$target), main="Density Plot: Stock Prices", ylab="
    Frequency", sub=paste("Skewness:", round(e1071::skewness(plotting_df$
    target), 2)))
91
92 ""
93
94
95
96 ""{ r}
97 sample_size = floor(0.8*nrow(stock_df))
98 set.seed(777)
99 train_stock_df = plotting_df%>%filter(Date<'2014-01-01')
100 test_stock_df = plotting_df%>%filter(Date>='2014-01-01')
101 ""
102
103
104
105
106

```

```

107
108
109 ““{ r}
110 linear_model <- lm(target ~ stock_val+log(1/Volume)+sentiment_score+sentiment_
    common, data=train_stock_df) # build linear regression model on full data
111 print(linear_model)
112 ““
113
114
115
116
117 ““{ r}
118 summary(linear_model)
119
120 ““
121
122
123 From the linear modelling we see that the most important factors are stock_val
    , Volume, and the sentiment_score of the company news. The common news
    data is nto considered to be incredibly informative to the Data.
124
125 ““{ r}
126 predicted <- predict(linear_model, test_stock_df)
127 actuals_preds <- data.frame(cbind(actuals=test_stock_df$target , predicted=
    predicted))
128 actuals_preds
129 ““
130
131 ““{ r}
132 correlation_accuracy <- cor(actuals_preds)
133
134 ““
135
136 ““{ r}
137 mae(actuals_preds_svr$actual , actuals_preds_svr$predicted)
138 ““
139
140 ““{ r}
141 library(class)
142 library(LiblineaR)
143 library(e1071)
144 SVR_model <- svm(target ~ stock_val+log(1/Volume)+sentiment_score+sentiment_
    common, data=train_stock_df) # build linear regression model on full data
145 print(SVR_model)
146 ““
147
148 ““{ r}
149 pred_svr <- predict(SVR_model, test_stock_df)
150 actuals_preds_svr <- data.frame(cbind(actuals=test_stock_df$target , predicted=
    =pred_svr))
151 actuals_preds_svr

```

```

152  “ “
153  “ “{ r }
154  correlation_accuracy_svr <- cor(actuals_preds_svr)
155
156  “ “

```

### 8.3.3 Neural Network

```

1  ———
2  title: "R Notebook"
3  output: pdf_output
4  ———
5
6
7  “ “{ r loading_data }
8  stock_df <- read.csv("Final_data_with_reddit.csv")
9  stock_df <- stock_df[,3:11]
10 head(stock_df)
11 “ “
12
13
14  “ “{ r }
15 x <- stock_df$sentiment_score
16 h<-hist(x, breaks=2, col="red", xlab="Sentiment",
17       main="Histogram with Normal Curve")
18 xfit<-seq(min(x),max(x),length=40)
19 yfit<-dnorm(xfit,mean=mean(x),sd=sd(x))
20 yfit <- yfit*diff(h$mids[1:2])*length(x)
21 lines(xfit, yfit, col="blue", lwd=2)
22 “ “
23
24  “ “{ r }
25 stock_df <- read.csv("Final_data_na_filtered")
26 stock_df <- stock_df[,3:11]
27 head(stock_df)
28 “ “
29
30
31  “ “{ r }
32 x <- stock_df$sentiment_score
33 h<-hist(x, breaks=10, col="red", xlab="Miles Per Gallon",
34       main="Histogram with Normal Curve")
35 xfit<-seq(min(x),max(x),length=40)
36 yfit<-dnorm(xfit,mean=mean(x),sd=sd(x))
37 yfit <- yfit*diff(h$mids[1:2])*length(x)
38 lines(xfit, yfit, col="blue", lwd=2)
39 “ “
40
41
42  “ “{ r }

```

```

43 library(corrplot)
44 M <- cor(stock_df[,4:8])
45 corrplot(M, method="circle")
46 ""
47
48 ""{ r}
49 stock_df <- stock_df%>%mutate(stock_val = (High+Low)/2)
50 ""
51
52
53
54
55
56
57
58
59
60 ""{ r}
61 library(keras)
62 stock_df <- stock_df%>%select(Date,company_name,company_name_full,stock_val,
    Volume,sentiment_score,sentiment_common,target)
63 stock_df <- stock_df%>%filter(target!=0)
64 stock_df
65 ""
66
67 ""{ r}
68 library(corrplot)
69 M <- cor(stock_df[,4:8])
70 corrplot(M, method="circle")
71 ""
72
73
74
75
76
77 ""{ r}
78 sample_size = floor(0.8*nrow(stock_df))
79 set.seed(777)
80
81 # randomly split data in r
82 picked = sample(seq_len(nrow(stock_df)),size = sample_size)
83 train_stock_df =stock_df[picked,]
84 test_stock_df =stock_df[-picked,]
85
86 train_stock <- train_stock_df%>%select(stock_val,Volume,sentiment_score,
    sentiment_common)
87 test_stock <- test_stock_df%>%select(stock_val,Volume,sentiment_score,
    sentiment_common)
88 train_target <- train_stock_df%>%select(target)
89 test_target <- test_stock_df%>%select(target)
90

```



```

91  “ “
92
93
94
95  “ “ { r }
96  train_data <- scale(train_stock)
97  col_means_train <- attr(train_data, "scaled:center")
98  col_stddevs_train <- attr(train_data, "scaled:scale")
99  test_data <- scale(test_stock, center = col_means_train, scale = col_stddevs_
    train)
100 “ “
101
102 “ “ { r }
103 build_model <- function() {
104
105   model <- keras_model_sequential() %>%
106     layer_dense(units = 64, activation = "relu",
107                 input_shape = dim(train_data)[2]) %>%
108     layer_dense(units = 64, activation = "relu") %>%
109     layer_dense(units = 1)
110
111   model %>% compile(
112     loss = "mse",
113     optimizer = optimizer_rmsprop(),
114     metrics = list("mean_absolute_error")
115   )
116
117   model
118 }
119
120 model <- build_model()
121 model %>% summary()
122
123 “ “
124
125
126 “ “ { r }
127 train_target <- data.matrix(train_target)
128 print_dot_callback <- callback_lambda(
129   on_epoch_end = function(epoch, logs) {
130     if (epoch %% 80 == 0) cat("\n")
131     cat(".")
132   }
133 )
134
135 epochs <- 300
136
137 # Fit the model and store training stats
138 history <- model %>% fit(
139   train_data,
140   train_target,

```

```

141 epochs = epochs ,
142 validation_split = 0.2 ,
143 verbose = 0 ,
144 callbacks = list(print_dot_callback)
145 )
146
147 “““
148
149
150
151
152 ““{ r}
153 library(ggplot2)
154
155 plot(history , metrics = "mean_absolute_error" , smooth = FALSE) +
156   coord_cartesian(ylim = c(0 , 5))
157 “““
158
159 ““{ r}
160 test_target <- data.matrix(test_target)
161 c(loss , mae) %<-%( model %>% evaluate(test_data , test_target , verbose = 0))
162
163 paste0("Mean absolute error on validation set: $" , sprintf("%.2f" , mae ))
164 “““
165
166 ““{ r}
167 test_target <- data.matrix(test_target)
168 c(loss , mae) %<-%( model %>% evaluate(test_data , test_target , verbose = 0))
169
170 paste0("Mean absolute error on test set: $" , sprintf("%.2f" , mae ))
171 “““
172
173 ““{ r}
174 test_predictions <- model %>% predict(test_data)
175 a <- as.data.frame(x = test_predictions[ , 1])
176 results <- a%>%mutate(actual = test_target)%>%mutate(prediction = test_
  predictions[ , 1]) %>% mutate(accuracy=100 - abs(((actual-prediction)/
  actual)*100))%>% select(actual , prediction , accuracy)
177 results <- results%>%filter(actual !=0)
178
179 “““
180
181 ““{ r}
182 results%>%filter(accuracy <90)
183 “““
184
185 ““{ r}
186 summary(results$accuracy)
187 “““

```

### 8.3.4 Recurrent Neural Network - LSTM

```

1  ———
2  title: "R Notebook"
3  output: pdf_output
4  ———
5
6  ““{ r}
7  library(keras)
8  library(dplyr)
9  ““
10
11 ““{ r loading_data}
12 stock_df <- read.csv("Final_data_with_reddit.csv")
13 stock_df <- stock_df[,3:11]
14 head(stock_df)
15 ““
16
17
18 ““{ r}
19 x <- stock_df$sentiment_score
20 h<-hist(x, breaks=2, col="red", xlab="Sentiment",
21       main="Histogram with Normal Curve")
22 xfit<-seq(min(x),max(x),length=40)
23 yfit<-dnorm(xfit,mean=mean(x),sd=sd(x))
24 yfit <- yfit*diff(h$mids[1:2])*length(x)
25 lines(xfit, yfit, col="blue", lwd=2)
26 ““
27
28
29 ““{ r}
30 library(corrplot)
31 M <- cor(stock_df[,4:8])
32 corrplot(M, method="circle")
33 ““
34
35 ““{ r}
36 stock_df <- stock_df%>%mutate(stock_val = (High+Low)/2)
37 ““
38
39
40
41
42 ““{ r}
43 library(keras)
44 stock_df <- stock_df%>%select(Date,company_name,company_name_full,stock_val,
45                               Volume,sentiment_score,sentiment_common,target)
46 stock_df <- stock_df%>%filter(target!=0)
47 ““
48 ““{ r}

```

```

49 stock_df$Date <- as.Date(stock_df$Date)
50 ""
51
52
53
54 ""{ r}
55 library(corrplot)
56 M <- cor(stock_df[,4:8])
57 corrplot(M, method="circle")
58 ""
59
60
61
62
63
64 ""{ r}
65 sample_size = floor(0.8*nrow(stock_df))
66 set.seed(777)
67 train_stock_df = stock_df%>%filter(Date<'2014-01-01')
68 test_stock_df = stock_df%>%filter(Date>='2014-01-01')
69 ""
70
71
72 ""{ r}
73 train_stock <- train_stock_df%>%select(stock_val, Volume, sentiment_score,
74   sentiment_common)
75 test_stock <- test_stock_df%>%select(stock_val, Volume, sentiment_score,
76   sentiment_common)
77 train_target <- train_stock_df%>%select(target)
78 test_target <- test_stock_df%>%select(target)
79 ""
80 ""{ r}
81 train_data <- scale(train_stock)
82 col_means_train <- attr(train_data, "scaled:center")
83 col_stddevs_train <- attr(train_data, "scaled:scale")
84 test_data <- scale(test_stock, center = col_means_train, scale =
85   col_stddevs_train)
86 ""
87
88 ""{ r}
89 train_X <- as.matrix(train_data)
90 train_y <- as.matrix(train_target)
91 test_X <- as.matrix(test_data)
92 test_y <- as.matrix(test_target)
93
94
95 dim(train_X) <- c(dim(train_X)[1],1,dim(train_X)[2])
96 dim(test_X) <- c(dim(test_X)[1],1,dim(test_X)[2])

```

```

97  “ “
98
99  “ “ { r }
100 library (" keras ")
101
102 rmsprop = optimizer_rmsprop ( lr = 0.00001, rho = 0.9, epsilon = 1e-08 )
103 model <- keras_model_sequential ()
104
105 model %>%
106   layer_lstm ( units           = 100,
107               input_shape      = c ( dim ( train_X ) [ 2 ], dim ( train_X ) [ 3 ] ),
108               return_sequences = TRUE,
109               ) %>%
110   layer_dropout ( 0.5 ) %>%
111   layer_lstm ( units           = 75,
112               return_sequences = FALSE,
113               ) %>%
114   layer_dropout ( 0.5 ) %>%
115   layer_flatten () %>%
116   layer_dense ( units = 64, activation = " relu " ) %>%
117   layer_dropout ( 0.5 ) %>%
118   layer_dense ( units = 1 )
119
120 model %>%
121   compile ( loss = ' mae ', optimizer = rmsprop )
122
123 model
124 “ “
125
126 “ “ { r }
127
128 epochs = 50
129
130
131 history <- model %>% fit (
132   train_X ,
133   train_y ,
134   epochs = epochs ,
135   batch_size = 10,
136   validation_split = 0.2,
137   verbose = 1,
138   shuffle = FALSE
139 )
140
141
142 “ “
143
144 “ “ { r }
145 model %>% save_model_hdf5 ( " model_lstm " )
146
147 “ “

```

```

148
149
150
151   “{ r}
152 pred_out <- model %>%
153   predict(test_X, batch_size = 10) %>%
154   .[,1]
155
156 length(pred_out)
157   “
158   “{ r}
159 plot(pred_out, metrics = "mean_absolute_error", smooth = FALSE) +
160   coord_cartesian(ylim = c(0, 5))
161
162 model
163   “
164
165   “{ r}
166 #predicted <- as.list(pred_out)
167 actual <- as.list(test_y)
168 a <- as.data.frame(test_y)
169 tbl_1 <- a %>%
170   mutate(predicted = pred_out)
171 results <- tbl_1 %>% mutate(accuracy=100 - abs(((target-predicted)/target)*100)
172   ) %>% select(target, predicted, accuracy)
173 results <- results %>% filter(actual != 0)
174   “
175   “{ r}
176 summary(results$accuracy)
177   “
178
179   “{ r}
180 mae(results$target, results$predicted)
181   “
182
183   “{ r}
184 output_df <- test_stock_df %>% mutate(predicted = pred_out)
185   “
186
187   “{ r}
188 library(plyr)
189 accuracy_df <- output_df %>% mutate(accuracy=100 - abs(((target-predicted)/
190   target)*100))
191 accuracy_df = accuracy_df %>% filter(accuracy != min(accuracy))
192 company_accuracy <- ddply(accuracy_df, .(company_name), summarize,
193   group_accuracy=mean(accuracy))
194   “
195

```

```

196  ““{ r}
197  company_accuracy%>%
198    ggplot() +
199    geom_bar(aes(x = company_name, y = group_accuracy), color = "darkred", stat =
      "identity")+
200    xlab('Company Name') +
201    ylab('Mean Accuracy')+ theme(axis.text.x = element_text(angle = 90, hjust =
      1))
202  ““
203
204  ““{ r}
205  IBM <- output_df%>%filter(company_name=='IBM')
206  Apple <- output_df%>%filter(company_name=='AAPL')
207  Walgreen <- output_df%>%filter(company_name=='WBA')
208  ““
209
210
211  ““{ r}
212  library(ggplot2)
213  library(dplyr)
214  library(plotly)
215  library(hrbrthemes)
216  output_df$Date <- as.Date(output_df$Date)
217
218  # Usual area chart
219  p <- Walgreen %>%
220    ggplot() +
221    geom_line(aes(x = Date, y = target), color = "darkred")+
222    geom_line(aes(x = Date, y = predicted), color="steelblue", linetype="twodash
      ") +
223    xlab('Dates') +
224    ylab('Stock Price—Walgreen')
225  p
226
227
228  ““
229
230  ““{ r}
231  library(ggplot2)
232  library(dplyr)
233  library(plotly)
234  library(hrbrthemes)
235  output_df$Date <- as.Date(output_df$Date)
236
237  # Usual area chart
238  p <- IBM %>%
239    ggplot() +
240    geom_line(aes(x = Date, y = target), color = "darkred")+
241    geom_line(aes(x = Date, y = predicted), color="steelblue", linetype="twodash
      ") +
242    xlab('Dates') +

```

```

243   ylab('Stock Price-IBM')
244 p
245
246 ' ' '
247
248 ' ' '{ r }
249 library(ggplot2)
250 library(dplyr)
251 library(plotly)
252 library(hrbrthemes)
253 output_df$Date <- as.Date(output_df$Date)
254
255 # Usual area chart
256 p <- Apple %>%
257   ggplot() +
258     geom_line(aes(x = Date, y = target), color = "darkred", alpha = 0.5)+
259     geom_line(aes(x = Date, y = predicted), color="steelblue", linetype="twodash") +
260     xlab('Dates') +
261     ylab('Stock Price-Apple')
262 p
263 ' ' '
264
265 ' ' '{ r }
266 library(ggplot2)
267 library(dplyr)
268 library(plotly)
269 library(hrbrthemes)
270 output_df$Date <- as.Date(output_df$Date)
271
272 png(filename="stock_plot_IBM.png", width=2500, height=500)
273 # Usual area chart
274 p <- IBM %>%
275   ggplot() +
276     geom_line(aes(x = Date, y = target, group = ), color = "red", alpha = 0.5)+
277     geom_line(aes(x = Date, y = predicted), color = "black") +
278     xlab('Dates') +
279     ylab('Stock Price')
280 p
281
282 ggsave(file="stock_plot_IBM.png", width=25, height=10, dpi=300)
283 ' ' '
284
285 ' ' '{ r }
286 q <- ggplotly(p)
287 q
288
289 ' ' '

```



## 9 REFERENCES

- [1] Anshul Mittal and Arpit Goel. Stock Prediction Using Twitter Sentiment Analysis, 2011.
- [2] Wenyue Sun, Chuan Tian and Guang Yang. Network Analysis of the Stock Market, 2015.
- [3] Kaggle, Daily News for Stock Market Prediction.url: <https://www.kaggle.com/aaron7sun/stocknews>
- [4] NYTimes,The New York Times Developer Network.url: <https://developer.nytimes.com>
- [5] Predicting Stock Prices using LSTM.  
<https://towardsdatascience.com/predicting-stock-price-with-lstm-13af86a7494>
- [6] <https://keras.io>